

Integrating active sensing into reactive synthesis with temporal logic constraints under partial observations

Jie Fu¹ and Ufuk Topcu¹

Abstract—We introduce the notion of online reactive planning with sensing actions for systems with temporal logic constraints in partially observable and dynamic environments. With incomplete information on the dynamic environment, reactive controller synthesis amounts to solving a two-player game with partial observations, which has impractically computational complexity. To alleviate the high computational burden, online replanning via sensing actions avoids solving the strategy in the reactive system under partial observations. Instead, we only solve for a strategy that ensures a given temporal logic specification can be satisfied had the system have complete observations of its environment. Such a strategy is then transformed into one which makes control decisions based on the observed sequence of states (of the interacting system and its environment). When the system encounters a belief—a set including all possible hypotheses the system has for the current state—for which the observation-based strategy is undefined, a sequence of sensing actions are triggered, chosen by an active sensing strategy, to reduce the uncertainty in the system’s belief. We show that by alternating between the observation-based strategy and the active sensing strategy, under a mild technical assumption of the set of sensors in the system, the given temporal logic specification can be satisfied with probability 1.

Keywords: Reactive synthesis; Active sensing; Partial observation; Temporal logic.

I. INTRODUCTION

Control synthesis under partial observations has been an important topic since complete and precise information (about the system and environment states) during the execution of a controller is often not available in practice. However, synthesis methods for systems under partial observations are of high complexity and have limitations in their applications. With incomplete information, the problem of synthesizing a controller in a partially observable Markov decision process (POMDP) has been shown to be PSPACE-complete, even for finite planning horizons [9]. When the control specification is given in temporal logic and the environment is dynamic and possibly adversarial, the interaction between a system and its environment can be captured in a two-player partially observable game with infinite stages, for which the qualitative-analysis problem under finite-memory strategies is EXPTIME-complete [3].

For temporal logic constraints, synthesis algorithms for stochastic systems modeled as POMDPs have been studied

in [11], [12]. To deal with a partially observable, dynamic environment, synthesis algorithms for two-player game with partial observations have been developed under two qualitative correctness criteria [2], [4]: *sure-winning* and *almost-sure winning* controllers. A sure-winning controller ensures the satisfaction of a specification whereas an almost-sure winning controller is a randomized strategy and ensures satisfaction with probability 1. These solutions rely on a subset construction and has complexity exponential in the size of the state space [3], [5].

An interesting question that has not been investigated much is the following: Since the high computational complexity is caused by incomplete information, is it possible to reduce the computational effort and still ensure correctness of the control design by acquiring new information at run time? In this paper, we give a method that provides a partial, affirmative answer to this question. Particularly, we study a system with actions to obtain information, referred to as *sensing actions*, and show how to utilize these actions in a way that a given linear temporal logic (LTL) specification is satisfied almost surely with reduced computational effort.

The new approach in this paper is inspired by [10], where the authors propose a method of online planning with partial observations and sensing actions as a way to overcome such complexity since the system only needs to compute a strategy for a finite number of steps, and replans with new information obtained through sensing actions. For temporal logic specifications, online planning method in [10] has no correctness guarantee. We propose a similar framework of active sensing and reactive synthesis under temporal logic constraints. The basic approach is the following: During control execution, the system maintains a *belief*, which is a set of states it thinks the current state must be in based on its partial observation for the game history. The belief is updated under two cases: In one of these cases, the system or the environment makes a move, the belief is updated to the set of states possibly arrived at as a result of move. Alternatively, the system can activate a sensor, detecting the value of some propositional formula and revises its belief according to the additional information obtained through sensing. In the second case, the system applies an active sensing strategy. A sequence of sensor queries are made to obtain the most useful information for reducing the system’s uncertainty in the current state. The benefit of performing the combined active sensing and reactive planning is that we can indeed avoid solving a two-player zero-sum game with partial observations. Rather, we transform the sure-winning

This work is supported by AFOSR grant number FA9550-12-1-0302, ONR grant number N000141310778 and NSF CNS award number 1446479.

¹Jie Fu and Ufuk Topcu are with the Department of Electrical and Systems Engineering, University of Pennsylvania, Philadelphia, PA, 19104, USA jief, utopcu@seas.upenn.edu.

strategy for the system in the same game with *perfect observations*, into a *randomized, belief-based* strategy. By construction, the randomized strategy may not be defined for every belief the system can encounter at run time. During control execution, the system alternates between the randomized strategy and the active sensing strategy. We prove that if the set of available sensors meets a sufficient condition, the temporal logic specification can be satisfied with probability 1, i.e., almost surely.

The rest of the paper is organized as follows. We begin with some preliminaries and the formulation of the problem in section II. Section III presents the main results on synthesizing provably correct, online reactive controllers with sensing actions for temporal logic constraints. In Section IV we illustrate the method using a robot motion planning example in a partially observed environment.

II. PROBLEM FORMULATION AND PRELIMINARIES

A probability distribution on a finite set S is a function $D : S \rightarrow [0, 1]$ such that $\sum_{s \in S} D(s) = 1$. The set of probability distributions on a finite set S is denoted $\mathcal{D}(S)$. The support of D is the set $\text{Supp}(D) = \{s \in S \mid D(s) > 0\}$. Let Σ be a finite alphabet. Σ^* , Σ^ω , and Σ^+ are sets of strings over Σ with finite length, infinite length, and length greater than or equal 1, respectively. Given u and v in Σ^* , uv is the concatenation of u with v . A string $u \in \Sigma^*$ is a *prefix* of $w \in \Sigma^*$ (or $w \in \Sigma^\omega$) if there exists $v \in \Sigma^*$ (or $v \in \Sigma^\omega$) such that $w = uv$. For a string w , the set of symbols occurring infinitely often in w is denoted $\text{Inf}(w)$. The last symbol in a finite string w is denoted $\text{Last}(w)$.

A. Game, specification and strategies

Through abstraction for systems with continuous and discrete dynamics, the interaction of a system and its dynamic environment can be captured by a labeled finite-state transition system [7], [8]:

$$M = \langle S, \Sigma, \delta, s_0, \mathcal{AP}, L \rangle$$

where 1) $S = S_1 \cup S_2$ is the set of states. At each state in S_1 , the system takes an action. At each state in S_2 , the environment takes an action. 2) $\Sigma = \Sigma_1 \cup \Sigma_2$ is the set of actions. Σ_1 is the set of actions for the system, and Σ_2 is the set of actions for the environment. 3) s_0 is the initial state. 4) $\delta : S \times \Sigma \rightarrow S$ is the transition function. 5) $L : S \rightarrow 2^{\mathcal{AP}}$ is the labeling function that maps a state $s \in S$ to a set of atomic propositions $L(s) \subseteq \mathcal{AP}$ that evaluate true at s .

We use a fragment of LTL [1] to specify the desired system properties such as safety, reachability, liveness and stability. Given a temporal logic formula φ in this class, one can always represent it by a deterministic Büchi automaton (DBA) $\mathcal{A}_\varphi = \langle H, 2^{\mathcal{AP}}, \delta_\varphi, h_0, F_\varphi \rangle$ where H is the set of states, $2^{\mathcal{AP}}$ is the set of alphabet, $\delta_\varphi : H \times 2^{\mathcal{AP}} \rightarrow H$ is the transition function, h_0 is the initial state and F_φ is the set of final states. A word $w = a_0 a_1 \dots \in (2^{\mathcal{AP}})^\omega$ induces a state sequence $h_0 h_1 \dots \in H^\omega$ where $h_{i+1} = \delta_\varphi(h_i, a_i)$, for all $i \geq 0$. A word w is accepted in \mathcal{A}_φ if and only if the

state sequence $\rho \in H^\omega$ induced from w visits some states in F_φ infinitely often.

A product operation is applied to incorporate the temporal logic specification into the labeled transition system, giving rise to a two-player turn-based Büchi game between the system (player 1) and its environment (player 2):

$$G = \langle Q, \Sigma, T, q_0, F \rangle = M \times \mathcal{A}_\varphi$$

where the components are defined as follows.

- $Q = Q_1 \cup Q_2$ is the set of states, where $Q_1 = S_1 \times H$ and $Q_2 = S_2 \times H$.
- $T : Q \times \Sigma \rightarrow Q$ is the transition function. Given $(s, h) \in Q$, $\sigma \in \Sigma$, if $\delta(s, \sigma) = s'$, then $T(q, \sigma) = q'$ where $q' = (s', \delta_\varphi(h, L(s')))$.
- $q_0 = (s_0, \delta_\varphi(h_0, L(s_0)))$ is the initial state.
- $F \subseteq Q \times F_\varphi$ is a subset of states that determines a Büchi winning condition.

A *play* in G is either a finite sequence of interleaving states and actions $\rho = q_0 a_0 q_1 a_1 \dots q_n \in (Q \cup \Sigma)^* Q$ or an infinite sequence $\rho = q_0 a_0 q_1 a_1 \dots \in (Q \cup \Sigma)^\omega$ such that q_0 is the initial state and $T(q_i, a_i) = q_{i+1}$ for all $i \geq 0$. If ρ is finite, the last element of ρ is a state, denoted $\text{Last}(\rho)$. An infinite *play* ρ is *winning* for player 1 in G if and only if $\text{Inf}(\rho) \cap F \neq \emptyset$.

In game G , each state in Q is associated with a truth assignment to a set \mathcal{P} of predicates. Note that \mathcal{P} may not equal \mathcal{AP} . This association is captured by the *interpretation* function π such that for any $q \in Q$, for any predicate $p \in \mathcal{P}$, $\pi(q)(p) \in \{\text{true}, \text{false}\}$. We write $\pi(q) = \bigwedge_{p \in \mathcal{P}} \ell_p$ where $\ell_p = p$ if $\pi(q)(p) = \text{true}$ and $\ell_p = \neg p$ if $\pi(q)(p) = \text{false}$, \wedge , \neg are the logical connectives for conjunction and negation, respectively. In the set \mathcal{P} , there is a predicate t indicating whose turn it is to play: If $t = 1$, then the system takes an action, otherwise the environment makes a move. It is assumed that the value of t is globally observable, which means, the system always knows whose turn it is to play.

We consider the case when the system has partial observation of values for the set \mathcal{P} of predicates. Following [4], this partial observation can be defined by an equivalence relation over the set of states, denoted $\mathcal{R} \subseteq Q \times Q$. Two states q and q' are *observation-equivalent*, that is, $(q, q') \in \mathcal{R}$, if both q and q' provide the same state information observable by the system, i.e., the value of $p \in \mathcal{P}$ is observable at q if and only if it is observable at q' , and $\pi(q)(p) = \pi(q')(p)$. We denote the *observations of states* for the system by $\mathcal{O} \subseteq 2^Q$, which is defined by the observation-equivalence classes. Clearly, \mathcal{O} is a partition of the state space. We define an *observation function* $\text{Obs} : Q \cup \Sigma \rightarrow \mathcal{O} \cup \Sigma_1 \cup \{-\}$ such that 1) $q \in \text{Obs}(q)$; 2) for every $q_1, q_2 \in \text{Obs}(q)$, $(q_1, q_2) \in \mathcal{R}$; 3) if $\sigma \in \Sigma_1$, $\text{Obs}(\sigma) = \sigma$; and 4) if $\sigma \in \Sigma_2$, $\text{Obs}(\sigma) = -$. The last two properties express that the system observes (knows) which action it performed but does not directly observe the action of the environment. The information received by the system on the environment's action is from the effect of that action, reflected in the observed arrived state.

The observation sequence of a play $\rho = q_0 a_0 q_1 \dots$ is a sequence $\text{Obs}(\rho) = \text{Obs}(q_0) \text{Obs}(a_0) \text{Obs}(q_1) \dots$. It is worth mentioning that two states $q = (s, h)$ and $q' = (s', h')$ can be observation-equivalent even if $h \neq h'$. Therefore, two observation-equivalent ρ and ρ' can differ in their state projections onto the set Q of states in the specification automaton \mathcal{A}_φ .

Let $\text{Pref}(G)$ denote the set of finite prefixes of all plays in G , each of which ends with a state in Q . For both players 1 and 2, a *deterministic* strategy for player i is a function $f_i : \text{Pref}(G) \rightarrow \Sigma_i$ and a *randomized* strategy is a function $f_i : \text{Pref}(G) \rightarrow \mathcal{D}(\Sigma_i)$. We say that player i follows strategy f_i if for any finite prefix $\rho \in \text{Pref}(G)$ at which f_i is defined, player i takes the action $f_i(\rho)$ if f_i is deterministic, or an action $\sigma \in \text{Supp}(f_i(\rho))$ with probability $f_i(\rho)(\sigma)$ if f_i is randomized. Since the system has partial information of the states, it can only execute an *observation-based* strategy f_1 , in the sense that if for any two prefixes ρ and $\rho' \in \text{Pref}(G)$, if $\text{Obs}(\rho) = \text{Obs}(\rho')$, then $f_1(\rho) = f_1(\rho')$. A strategy is *memoryless* if and only if $f_i(\rho) = f_i(\text{Last}(\rho))$. For Büchi game G with complete information, there exists a *deterministic, memoryless* winning strategy for one of the players.

B. Partial observation, belief and sensing actions

With partial observations, the system keeps track of the play in the game by maintaining and updating a set $B \subseteq Q$ of states, referred to as the *belief*, which is the set of states the system thinks the game can be in, given the observation history. In which follows, we show how the belief is obtained and updated. The set of beliefs in the game is denoted $\mathcal{B} \subseteq 2^Q$. We define a function $\alpha : \text{Pref}(G) \rightarrow \mathcal{B}$ that maps a prefix of g into a *belief* as follows: given a prefix $\rho = q_0 a_0 \dots q_n$, the belief of the system is $\alpha(\rho) = \{\text{Last}(\rho') \in Q \mid \rho' \in \text{Pref}(G) \text{ and } \text{Obs}(\rho') = \text{Obs}(\rho)\}$.

During the interaction with the environment, the system's belief is updated in two ways: (i) The system applies a control action, obtains a new observation of the arrived state, and updates its belief to one in which the current state could be. (ii) The environment takes some action. The system obtains an observation $o \in \mathcal{O}$ of the arrived state, and subsequently updates its belief that includes its hypothesis for the current state. Formally, this process is called *belief update*, which can be captured by the function

$$\text{Update} : \mathcal{B} \times (\Sigma_1 \cup \{-\}) \times \mathcal{O} \rightarrow \mathcal{B}, \quad (1)$$

It is reminded that the symbol “ $-$ ” is the observation for an action of the environment. Given a belief B , the system takes an action $a \in \Sigma_1$ and gets an observation $o \in \mathcal{O}$. Then it updates its belief to $B' = o \cap \text{Update}(B, a, o) = \{q' \mid \exists q \in B \text{ such that } T(q, a) = q'\}$. If it is the environment's turn, after the environment takes some action, the system gets an observation $o \in \mathcal{O}$ and then updates its current belief B to $B' = \text{Update}(B, -, o) = o \cap \{q' \mid \exists q \in B, \exists \sigma \in \Sigma_2 \text{ such that } T(q, \sigma) = q'\}$.

We distinguish a set Γ of *sensing actions* for the system and explain how the sensing actions affects the system's belief as follows.

Definition 1: Consider the set \mathcal{P} of atomic propositions and the set Γ of sensing actions. For each sensing action $a \in \Gamma$, there exists at least one propositional formula ϕ over \mathcal{P} such that after applying the sensing action a , the truth value of ϕ is known. Depending on the value of ϕ , the system can partition a belief B into two subsets, expressed by

$$\text{Knows}(\phi, a, B) := (B', B \setminus B'),$$

where B' is the set of states in which ϕ evaluates true and $B \setminus B'$ is the set of states in which ϕ evaluates false. Hence, if ϕ is true, the belief is revised to be B' , otherwise to be $B \setminus B'$.

To capture both global and local sensing capabilities, for a given state q , we denote $\Gamma_q \subseteq \Gamma$ to be a set of sensing actions *enabled* at q . The set of sensing actions enabled at a belief $B \subseteq Q$ is $\bigcap_{q \in B} \Gamma_q$.

The following assumption is made for sensing actions.

Assumption 1: A sensing action will not change the value of variables and/or predicates in \mathcal{P} .

The assumption is not restrictive because if an action introduces both physical and epistemic changes, we simply consider it as an ordinary control action and include it into Σ_1 . We call an action in Γ *sensing* to emphasize that it provides information of the current state, and an action in Σ *physical* to emphasize it changes the state of the game. We assume that at each turn of the system, it can either choose a physical action, or several sensing actions followed by a physical action.

We solve the following problem in this paper.

Problem 1: Given a two-player turn-based Büchi game $G = \langle Q, \Sigma, T, q_0, F \rangle$, and a set Γ of sensing actions, design an observation-based strategy $f : Q^* \rightarrow D(\Sigma_1) \cup \Gamma^*$ with which the specification is satisfied with probability 1, i.e., almost surely, whenever such a strategy exists.

III. MAIN RESULTS

For games with partial information, algorithms in [5] can be used to synthesize observation-based controllers which ensure given temporal logic specifications are satisfied surely, or almost surely, i.e., with probability 1, whenever such controllers exist. In this paper, we only consider the cases in which observation-based controllers do not exist and thus require additional information at run time for satisfying given temporal logic specifications. We distinguish two phases in the online planning: *Progress* phase and *sensing* phase. As the names suggest, during the progress phase, the system takes physical actions in order to satisfy the temporal logic constraints, and during the sensing phase, the system takes sensing actions to reduce the uncertainty in its belief for the current game state. The transition from one phase to another will be explained after we introduce the methods for synthesizing strategies used in both phases.

A. A belief-based strategy for making progress

For a game with partial observation, we aim to synthesize a belief-based, memoryless and randomized strategy $f_P : \mathcal{B} \rightarrow \mathcal{D}(\Sigma_1)$ that can be applied for making progress towards satisfying the given LTL fragment formula φ .

In the two-player Büchi game G , the deterministic sure-winning strategy $WS : Q \rightarrow \Sigma_1$ can be computed (with methods in [6]) but requires complete information to execute at run time. The belief-based strategy f_P is constructed from the sure-winning strategy WS in the following way: Let $Win_1 \subseteq Q$ be the set of states at which WS are defined. Given $B \in \mathcal{B}$, let

$$\begin{aligned} \text{Progress}(B) &= \bigcup_{q \in B} WS(q), \text{ and} \\ \text{allow}(B) &= \bigcap_{q \in B} \text{allow}(q), \\ \text{where } \text{allow}(q) &= \{\sigma \in \Sigma_1 \mid T(q, \sigma) \in Win_1\}. \end{aligned}$$

For each state $q \in B$, the sure-winning strategy will suggest action $WS(q)$ to be taken by the system, which is then included into a set $\text{Progress}(B)$. The set $\text{allow}(B)$ is a set of actions with the following property: No matter in which state of B the game is, by taking an action in $\text{allow}(B)$, the next state will still be one for which the sure-winning strategy is defined. Then, if $\text{Progress}(B) \subseteq \text{allow}(B)$, we let $f_P(B)(\sigma) = \frac{1}{|\text{Progress}(B)|}$ for each $\sigma \in \text{Progress}(B)$. Otherwise, f_P is undefined for B . Note that since the computation f_P can be essentially reduced to computing the interaction of two sets, there is no need to compute f_P for all possible subset of Q . Rather, we can efficiently compute f_P for each belief B encountered at run time.

We have transformed the sure-winning strategy with complete information in the Büchi game into a randomized, belief-based strategy. During control execution, the system maintains its current belief. At each turn of the system, after applying an action $\sigma \in \Sigma_1$ at the state B , the system receives an observation $o \in \mathcal{O}$, updates its belief to $B' = \text{Update}(B, \sigma, o)$. When it is a move made by the environment, the system obtains another observation $o' \in \mathcal{O}$, updates its belief to $B'' = \text{Update}(B', -, o')$. The system applies $f_P(B'')$ as long as f_P is defined for B'' . When f_P is undefined for the current belief B , then we switch to the sensing phase for actively acquiring more information to reduce the uncertainty in its current belief.

B. An active sensing strategy for reducing uncertainty

During the progress phase with the randomized, belief-based strategy f_P , if the system runs into a belief at which f_P is undefined, it needs to update its belief through sensing until either it finds itself in a state for which f_P is defined, or it cannot further refine its belief: A belief B cannot be refined if for any sensing action a enabled at B and for any formula ϕ such that $(B_1, B_2) = \mathbf{Knows}(\phi, a, B)$, it holds that for either $i = 1$ or $i = 2$, $B_i = B$. We represent the process of belief revision with sensing actions as a tree structure, referred to as a *belief revision tree*, and then

propose a synthesis method for an active sensing strategy using the belief revision tree.

Given a belief $B^o \in \mathcal{B}$, the *belief revision tree* with the root B^o is a tuple $\text{BRTree}(B^o) = \langle \mathcal{N}, \mathcal{E} \rangle$, where \mathcal{N} is the set of nodes in the tree, consisting a subset of beliefs, and $\mathcal{E} \subseteq \mathcal{N} \times \Gamma \times \mathcal{N}$ is the set of edges. It is constructed as follows.

- 1) The root of the tree is B^o .
- 2) At each node $B \in \mathcal{N}$, for each enabled sensing action $a \in \Gamma_B$, if there exists a formula ϕ such that $(B_1, B_2) = \mathbf{Knows}(\phi, a, B)$ and both B_1, B_2 are not empty, then we add two children B_1, B_2 of B , and include edges $(B, a, B_1), (B, a, B_2)$ into the edges \mathcal{E} .
- 3) A node B is a leaf of the tree if and only if either
 - 1) B cannot be further revised by any sensing action,
 - or 2) f_P is defined for B .

The active sensing strategy $f_S : \mathcal{B} \rightarrow \Gamma$ is computed as follows. First, in the tree $\text{BRTree}(B^o)$, we compute a set of target nodes $\text{Reach} \subset \mathcal{N}$ such that a node B' is included in Reach if and only if $f_P(B')$ is defined. The objective is to apply the least number of sensing actions in order to reach a belief in Reach for which f_P is defined. For this purpose, we have the following recursion:

- 1) $X_0 = \text{Reach}, i = 0$.
- 2) $X_{i+1} = X_i \cup \{B \in \mathcal{B} \mid \exists a \in \Gamma, \text{ such that } \forall B' \in \mathcal{B}, (B, a, B') \in \mathcal{E}, B' \in X_i\}$ and let $f_S(B) = a$. In other words, a belief B is included into X_{i+1} if there exists a sensing action a such that when a is applied at B , no matter which belief the system might reach, it must be in X_i .
- 3) Until i is increased to some number $m \in \mathbb{N}$ such that $X_{m+1} = X_m$, we output the sensing strategy f_S obtained so far.

We denote $X_m = \text{attr}(\text{Reach})$, following the notion of an *attractor* of the set Reach . For any state in $\text{attr}(\text{Reach})$, there exists a sensing strategy f_S such that for *whatever outcome* resulted by applying sensing actions, the system can arrive at some belief in Reach in finitely many steps by following f_S . Furthermore, it can be proven that f_S minimizes the number of sensing actions required for the sensing phase under the constraint that the system will not run into a dead end, which is a belief that cannot be further refined yet is undefined by f_P . The number of sensing actions during the sensing phase is upper bounded by the index i for which $B^o \in X_i$ and $B^o \notin X_{i-1}$. The proof follows from the property of attractor [6] and is omitted here.

Remark: It is worth mentioning that for a given belief B , the active sensing strategy is unique. Thus, we can store and continuously update a set of active sensing strategies synthesized at run time: When the system encounters a belief B for which f_P is undefined but it has seen before, it can use the stored active sensing strategy for B without recomputing a new one. For a large-scale system with a large number of sensing actions, one can also pre-compute a library of active sensing strategies and then augment the library with

new active sensing strategies computed at run time.

C. A composite, almost-sure winning strategy

At run time, the system alternates between strategy f_P for making progress and strategy f_S for refining its belief. We name the system's strategy at run time a *composite* strategy, denoted $f : \mathcal{B} \rightarrow D(\Sigma_1) \cup \Gamma$, defined by,

$$f(B) = \begin{cases} f_P(B) & \text{if } f_P(B) \text{ is defined.} \\ f_S(B) & \text{if } f_S(B) \text{ is defined.} \end{cases} \quad (2)$$

Note that by construction, the domains of f_P and f_S is always disjoint.

The following assumption provides a sufficient condition for avoiding dead-ends at run time.

Assumption 2: For each state B encountered during the progress phase, if $f_P(B)$ is undefined, then $f_S(B)$ is defined.

Since we cannot predict which beliefs the system might have during control execution with online planning, in the extreme case, for each predicate $p \in \mathcal{P}$, we need to have a sensing action or a combination of sensing actions to detect its truth value. However, this condition is not necessary and may include some sensing actions that will never be used at run time. As the system does not need to know the *exact* state by extensive sensing, it is at the system's disposal whether to apply a sensing action and what shall be applied.

Next we prove the correctness of the composite strategy. To this end, we recall some property in the solution for Büchi games with complete information from [6]: The winning region of the Büchi game G can be partitioned as $\text{Win}_1 = \bigcup_{i=0}^m W_i$ for some $m \in \mathbb{N}$, $m \geq 0$. For any state $q \in \text{Win}_1$, there exists a unique ordinal i such that $q \in W_i$. If $q \in Q_1 \cap W_i$ for some $0 < i \leq m$, then the winning strategy on q outputs $\sigma \in \Sigma_1$, with which the system reaches a state $q' \in W_{i-1} \cap Q_2$. If $i = 0$, then with the action $\text{WS}(q)$, we arrive at a state $q' \in \text{Win}_1$. If $q \in Q_2$, then for any action $\sigma \in \Sigma_2$ enabled at q , $T(q, \sigma) \in W_{i-1}$ if $i \neq 0$, or $q' \in \text{Win}_1$ otherwise.

Lemma 1: Given a game $G = \langle Q, \Sigma, T, q_0, F \rangle$. Let $B_0 = \text{Obs}(q_0)$ be the initial belief. If Assumption 2 is satisfied and $q_0 \in \text{Win}_1$, the composite strategy f defined by (2) ensures that some states in F of G is infinitely often visited with probability 1.

Proof: Consider an arbitrary belief $B \in \mathcal{B}$ for which f_P is defined. By definition of f_P , for each $\sigma \in \text{Progress}(B)$, the probability of choosing action σ is $\frac{1}{u}$, where $u = |\text{Progress}(B)|$. If the actual state is q and $q \in W_i$, for some $i \neq 0$, then with probability $\frac{1}{u}$, the system will reach a state in W_{i-1} . Thus, the probability of the next state being in W_{i-1} is $\frac{1}{u} \geq \frac{1}{|Q|} > 0$. For other $\sigma' \in f_P(B)$, $\sigma' \neq \text{WS}(q)$, the next state after taking σ' is in W_j for some $0 \leq j \leq m$. Let $\text{Pr}(q, \Diamond^i W_0)$ denote the probability of reaching W_0 from state q in i turns. When system applies the strategy f , it is $\text{Pr}(q, \Diamond^i W_0) \geq (\frac{1}{|Q|})^i > 0$ and the probability of *not* reaching W_0 in i turns is *less than or equal to* $1 - (\frac{1}{|Q|})^i \leq 1 - (\frac{1}{|Q|})^{m+1} = r < 1$ where $m+1$ is the total number of partitions in Win_1 . If after

i steps the state is not in W_0 , it must be in W_j for some $0 < j \leq m$, and again the probability of not reaching W_0 in m steps is less than or equal to r . Therefore, under the policy f , the probability *eventually* reaching W_0 from any state $q \in \text{Win}_1$ is $\text{Pr}(q, \Diamond W_0) = \lim_{k \rightarrow \infty} \text{Pr}(q, \Diamond^k W_0) = \lim_{k \rightarrow \infty} (1 - \text{Pr}(q, \neg \Diamond^k W_0)) = \lim_{k \rightarrow \infty} (1 - r^{k/m}) = 1 - \lim_{k \rightarrow \infty} r^{k/m} = 1$.

Once entering W_0 , the system will take an action to remain in Win_1 , and the above reasoning applies again. In this way, in the absence of dead ends (Assumption 2), the system can revisit the set W_0 of states with probability 1 by following the composite strategy f . Since $W_0 \subseteq F$, the probability of system always eventually visiting some states in F is 1. ■

To conclude this section, Algorithm 1 describes the procedure of online planning with sensing actions.

Algorithm 1: PlanningWithSensingActions

Input: A labeled finite-state transition system $M = \langle S, \Sigma, \delta, s_0, \mathcal{AP}, L \rangle$ and a DBA $\mathcal{A}_\varphi = \langle H, 2^{\mathcal{AP}}, \delta_\varphi, h_0, F_\varphi \rangle$ representing a LTL fragment formula. An observation function $\text{Obs} : Q \cup \Sigma \rightarrow \mathcal{O} \cup \Sigma_1 \cup \{-\}$.

Output: An online planning strategy that ensures φ is satisfied with probability 1.

```

begin
   $G = \langle Q, \Sigma, T, q_0, F \rangle := M \times \mathcal{A}_\varphi$ ;
   $(\text{WS}, \text{Win}_1) := \text{GameSolve}(G)$  /* Solve  $G$ 
    with perfect information */
   $B_0 \leftarrow \text{Obs}(q_0)$  /* The initial belief */
   $B \leftarrow B_0$ 
  while True do
    if  $B$  is a turn for the system then
      if  $\text{Progress}(B) \subseteq \text{allow}(B)$  then
         $\sigma \leftarrow \text{Random.Choice}(\text{Progress}(B))$ ;
        Receive a new observation  $o$ ;
         $B \leftarrow \text{Update}(B, \sigma, o)$ ;
      else
         $f_s := \text{GetSensingStrategy}(B, \Gamma)$ ;
         $\sigma \leftarrow f_s(B)$  /* Sensing action  $\sigma$ 
          with effect  $\text{Knows}(\phi, \sigma, B)$  */
         $(B_1, B_2) = \text{Knows}(\phi, \sigma, B)$ 
        if  $\phi = \text{true}$  then  $B \leftarrow B_1$ ;
        else  $B \leftarrow B_2$ ;
      else
        Environment makes a move and the system
        observes  $o$ ;  $B \leftarrow \text{Update}(B, -, o)$ 

```

Fig. 1: Algorithm: PlanningWithSensingActions

IV. EXAMPLES

We apply the algorithm to a robotic motion planning example, which is a variant of the so-called “Wumpus game” in a 7×7 gridworld. Figure 2 consists of one mobile robot, one monster called “Wumpus”. The robot is capable of moving in eight compass directions with actions ‘N’, ‘S’, ‘E’, ‘W’, ‘NE’, ‘NW’, ‘SE’, ‘SW’ (horizontally, vertically and diagonally), one step at a time. The robot and the Wumpus does not move concurrently. The Wumpus can move in four compass directions with actions ‘N’, ‘S’, ‘E’ and ‘W’ within a restricted area Region and emits stench to its surrounding cells. The objective of the robot is to infinitely revisit region R_1 , R_2 , and R_3 in this order, while avoiding running into

the Wumpus. Formally, the temporal logic formula is $\varphi = \Box \Diamond (x_r, y_r) = R_1 \wedge \Diamond ((x_r, y_r) = R_2 \wedge \Diamond (x_r, y_r) = R_3) \wedge \Box \neg (x_r = x_w \wedge y_r = y_w)$ where $(x_r, y_r), (x_w, y_w)$ are the positions of the robot and the Wumpus, respectively. Yet, the robot only knows his own position. For this case of partial observation, without the inclusion of sensing actions, it can be shown that with the algorithms in [5], observation-based, sure-winning strategies and almost-sure winning strategies do not exist.

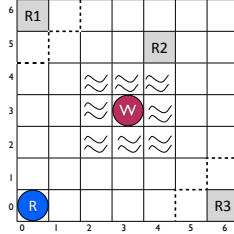


Fig. 2: The gridworld with a robot (R) and the Wumpus (W). The grey cells are regions R_1 , R_2 and R_3 . The Wumpus is restricted to the area inside the dash lines. The stenches emitted by the Wumpus are represented by waves.

Here, we introduce a set of sensing actions to the game. For the robot to know the position of the moving obstacles, it needs to apply a sensing action — $\text{smell}(x, y)$ to detect if there exists stench at cell (x, y) . Thus, when the robot applies $\text{smell}(x, y)$, if the result is True, then the Wumpus must be some cells in the set $S = \{(x', y') \mid x' \leq x+1, y' \leq y+1, x', y' \in \mathbb{N}\} \cap \text{Region}$. Otherwise, it is not possible that the Wumpus is in any cell in S .

We illustrate how the robot updates his belief using sensing action $\text{smell}(x, y)$ where (x, y) is a cell in the gridworld. Suppose that the robot does not know where the Wumpus is and hypothesizes it can be in any cell in the Region. Once it applies the sensing action (2,2), since the cell has stench and the sensor returns True. Then, immediately the robot will know the Wumpus is in one of the cells in the set $S = \{(1, 1), (2, 1), (3, 1), (1, 2), (2, 2), (3, 2), (3, 1), (3, 2), (3, 3)\}$, because only if the Wumpus is in a cell of S , there can be stench in cell (2,2).

From the numerical experimental result, after 1000 steps (a step includes either a robot's (sensing or physical) action or a movement of the Wumpus), the robot visited the set F in the formulated two-player game G 14 times and can continue to visit F infinite often. In Figure 3 we show the belief updates by applying alternatively the exploitation strategy and active sensing strategy for the initial 100 steps. It is observed that the maximum cardinality of the belief set is 43 over the control execution, which means that the robot thinks the Wumpus can be in any cell in its restricted region. However, if there is no danger of running into the Wumpus in a few next steps, there is no need to exercising any sensing action. The implementations are in Python on a desktop with Intel(R) Core(TM) i5 processor and 16 GB of

memory. The average time for the robot making a decision is 8.55×10^{-4} seconds. The computation of the product game took 40.14 seconds and the winning strategy under complete information is computed within 14 seconds.

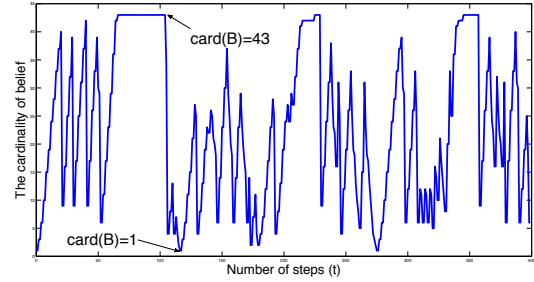


Fig. 3: The update in the number of possible Wumpus' positions in the system's belief.

V. CONCLUSIONS

Our work shows that when additional information can be obtained through sensing actions, one can transform a sure-winning strategy with complete information to a belief-based, randomized strategy, which is then combined, at run time, with an active sensing strategy to ensure a given temporal logic specification is satisfied with probability 1. The synthesis method avoids a subset construction for solving games with partial information. Meanwhile, the active sensing strategy leads to a cost-efficient way of sensor design: Although we require a sufficient set of sensing actions to avoid dead-ends at run time, the system minimizes the usage of sensing actions by asking the most revealing queries, depending on what specification is to be satisfied, and how much uncertainty the system has about the game state at run time. In future work, we will consider more examples for practical robotic motion planning under partial observations. It is also important to consider the uncertainty in the sensors. For example, a sensor query might return a probabilistic distribution over a set of states, rather than a

binary answer to proposition logical formulae considered herein. For this extension, we are currently investigating modifications that need to be made to account for delays, uncertainty in the information provided by the sensors.

REFERENCES

- [1] Rajeev Alur and Salvatore La Torre. Deterministic generators and games for LTL fragments. *ACM Transactions on Computational Logic*, 5(1):1–25, January 2004.
- [2] A Arnold, A Vincent, and I Walukiewicz. Games for synthesis of controllers with partial observation. *Theoretical Computer Science*, 303(1):7–34, 2003.
- [3] Krishnendu Chatterjee and Laurent Doyen. The complexity of partial-observation parity games. In *Logic for Programming, Artificial Intelligence, and Reasoning*, pages 1–14. Springer, 2010.
- [4] Krishnendu Chatterjee and Laurent Doyen. Partial-Observation Stochastic Games: How to Win When Belief Fails. *Annual IEEE Symposium on Logic in Computer Science*, pages 175–184, June 2012.
- [5] Krishnendu Chatterjee, Laurent Doyen, Thomas A Henzinger, and Jean-François Raskin. Algorithms for omega-regular games with imperfect information. *Logical Methods in Computer Science*, 3(4):1–23, 2007.
- [6] Erich Grädel, Wolfgang Thomas, and Thomas Wilke, editors. *Automata Logics, and Infinite Games: A Guide to Current Research*. Springer-Verlag New York, Inc., New York, NY, USA, 2002.
- [7] M. Kloetzer and C. Belta. A fully automated framework for control of linear systems from temporal logic specifications. *Automatic Control, IEEE Transactions on*, 53(1):287–297, Feb 2008.
- [8] Hadas Kress-Gazit, Tichakorn Wongpiromsarn, and Ufuk Topcu. Correct, reactive robot control from abstraction and temporal logic specifications. *IEEE Robotics and Automation Magazine*, 18:65–74, 2011.
- [9] Michael Lederman Littman. *Algorithms for sequential decision making*. PhD thesis, Brown University, 1996.
- [10] Guy Shani and Ronen I Brafman. Replanning in domains with partial information and sensing actions. In *IJCAI*, volume 2011, pages 2021–2026, 2011.
- [11] Rangoli Sharan. *Formal methods for control synthesis in partially observed environments : application to autonomous robotic manipulation. Dissertation (Ph.D.)*, California Institute of Technology. PhD thesis, California Institute of Technology, 2014.
- [12] Tichakorn Wongpiromsarn and Emilio Frazzoli. Control of probabilistic systems under dynamic, partially known environments with temporal logic specifications. In *Proceedings of the 51th IEEE Conference on Decision and Control, CDC 2012, December 10-13, 2012, Maui, HI, USA*, pages 7644–7651. IEEE, 2012.